

Similarity Search in Large Collections of Biometric Data

Pavel Zezula – Michal Batko – Vlastislav Dohnal – David Novak – Jan Sedmidubsky

Faculty of Informatics, Masaryk University

Botanicka 68a, 602 00 Brno

Czech Republic

[zezula, batko, dohnal, xnovak8, xsedmid]@fi.muni.cz

ABSTRACT

The new field of terrorism informatics is not only multidisciplinary but also requires new approaches to processing its complex underlying collections of data. In this paper, we introduce Multi-Feature Indexing Network (MUFIN) which is able to manage and search a large class of digital data according to the metric model of similarity. It is also scalable with respect to the data volume as well as the query execution throughput. At the same time, its performance can be tuned by mapping the search structure on a needed computation infrastructure. We illustrate the capabilities of MUFIN by outlining some of its current applications. Finally, we show how MUFIN can be used to deal with many biometric data and illustrate such possibility on the face retrieval application.

1.0 INTRODUCTION

The multidisciplinary field of terrorism informatics – aiming at accurately and efficiently monitoring, analyzing, predicting and preventing terrorist activities – has recently experienced tremendous growth. Accordingly, the development and use of advanced information technologies, including methodologies, models and algorithms, infrastructure, systems, and tools for national/international and homeland security related applications have provided promising new directions for study. Terrorism informatics is a highly interdisciplinary and comprehensive field. The wide variety of methods used in terrorism informatics are derived from computer science, informatics, statistics, mathematics, linguistics, and social sciences, and these methods must manage collections of huge amounts of many types of digital information from varied and multiple sources. Information fusion and information technology analysis techniques – which include data mining, data integration, language translation technologies, and image and video processing – require suitable search mechanisms to play central roles in the prevention, detection and remediation of terrorism.

An important category of data for terrorism informatics concerns the biometric data. Biometrics refers to the use of a person's physical characteristics or personal traits – such as fingerprints, faces, voices, or handwritten signatures – for identification. Nowadays, there is a technology capable of producing large amount of digital biometric signatures, and it is believed that biometric-based systems will become increasingly important tools for identifying known and suspected terrorists. Biometric-based systems could provide automatic, nearly instantaneous identification of a person by converting the biometric into a digital form and then comparing it against a computerized database. They can be used to improve security and thereby help safeguard our communities against future terrorist attacks in critical areas such as: (1) controlling access to sensitive facilities at airports, (2) preventing identity theft and fraud in the use of travel documents, and (3) identifying known or suspected terrorists. There is no high-tech silver bullet to solve the terrorism problem, and several biometric characteristics have to be typically combined to get reliable results. Exact match in biometric collections have very little meaning and only a relative ordering of database objects with respect to a reference can be achieved by means of a ranking function or a similarity measure. Consequently, standard search technology fails and new similarity search engines must be developed. In other words, the central issue determining effectiveness of such technologies is the way they deal with similarity of processed entities.

Report Documentation Page				Form Approved OMB No. 0704-0188	
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE OCT 2009		2. REPORT TYPE N/A		3. DATES COVERED -	
4. TITLE AND SUBTITLE Similarity Search in Large Collections of Biometric Data				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Faculty of Informatics, Masaryk University Botanicka 68a, 602 00 Brno Czech Republic				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release, distribution unlimited					
13. SUPPLEMENTARY NOTES See also ADA562563. RTO-MP-MSG-069 Current Uses of M&S Covering Support to Operations, Human Behaviour Representation, Irregular Warfare, Defence against Terrorism and Coalition Tactical Force Integration (Utilisation actuelle M&S couvrant le soutien aux operations, la representation du comportement humain, la guerre asymetrique, la defense contre le terrorisme et l'integration d'une force tactique de coalition). Proceedings of the NATO RTO Modelling and Simulation Group Symposium held in Brussels, Belgium on 15 and 16 October 2009., The original document contains color images.					
14. ABSTRACT The new field of terrorism informatics is not only multidisciplinary but also requires new approaches to processing its complex underlying collections of data. In this paper, we introduce Multi-Feature Indexing Network (MUFIN) which is able to manage and search a large class of digital data according to the metric model of similarity. It is also scalable with respect to the data volume as well as the query execution throughput. At the same time, its performance can be tuned by mapping the search structure on a needed computation infrastructure. We illustrate the capabilities of MUFIN by outlining some of its current applications. Finally, we show how MUFIN can be used to deal with many biometric data and illustrate such possibility on the face retrieval application.					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT SAR	18. NUMBER OF PAGES 14	19a. NAME OF RESPONSIBLE PERSON
a. REPORT unclassified	b. ABSTRACT unclassified	c. THIS PAGE unclassified			

As the biometric systems are trying to substitute people in identification or authentication processes they must deal with similarity because the ability to access similarity lies close to the core of human cognition. In fact, the sense of sameness is the very keel and backbone of our thinking. Similarity serves an organizing principle by which individuals classify objects, form concepts, and make categorizations. In practice, people organize, group and categorize things based on their degree of similarity and separate them based on their degree of difference or dissimilarity. The digital biometric characterization of people must be processed by analogy.

Though similarity has been the focus of many investigations in psychology, cognitive and behavioral sciences for over 100 years and it is one of the most important and researched constructs, the application of these concepts to audio-visual and biometric digital data is complicated mainly for the following three reasons. First, we need a single system which is able to process many, possibly all, forms of similarity which occur in the range of given application, because the cost of building such system is excessive. Second, the system must be able to provide performance required by the given application – online processing is typically needed. Third, the properties of the system must scale with respect to the volume of data as well as the user processing workload.

In this paper, we present basic concepts and implementation strategies of a similarity search engine called MUFIN (Multi-Feature Indexing Network). It is an extensible, scalable, and infrastructure independent engine, based on the metric space vision of similarity. We describe and explain its architecture in Section 2. Section 3 is devoted to the implementation background of the MUFIN system. The properties of extensibility and scalability are demonstrated in Section 4 by running applications operating on several collections of image data. Section 5 starts by presentation of a face retrieval application and demonstrates how many other biometric applications can be developed by analogy. The paper concludes in Section 6.

2.0 MUFIN ARCHITECTURE

From a general point of view, the search problem has three dimensions shown in Figure 1: (1) data and query types, (2) index structures and search algorithms, and (3) infrastructure to run the system on. The MUFIN offers a complex solution that is highly flexible in all three aspects.

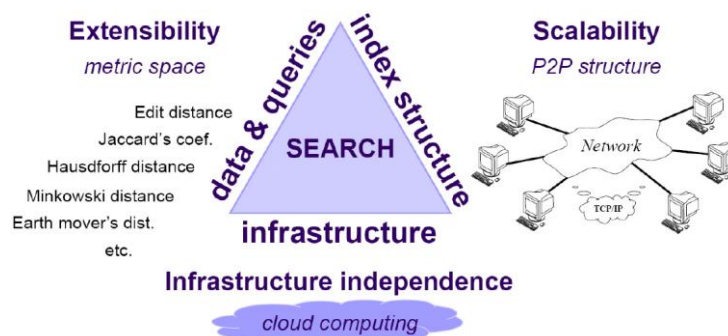


Figure 1: The three aspects of the search problem in the MUFIN system.

The data processed in MUFIN are modeled using the *metric space* approach. It means that the data can be practically anything that has a digital representation if a method to measure their similarity exists. More formally, the mathematical metric space is a pair (D, d) , where D refers to a domain of objects and d is a function able to compute the distance between any pair of objects from D . It is assumed that the smaller the distance, the closer or more similar the objects are and the zero distance means that the objects are identical. For any three distinct objects $x, y, z \in D$, the distance must satisfy properties of reflexivity, $d(x, x) = 0$, strict positivity, $d(x, y) > 0$, symmetry, $d(x, y) = d(y, x)$, and triangle inequality, $d(x, y) \leq d(x, z) + d(z, y)$.

+ $d(z, y)$. Such definition allows us to specify a number of similarity queries, including the similarity range and nearest neighbor queries. Many similarity functions defined for various data types satisfy these three natural properties. For example, the commonly used Euclidean distance on 2D or 3D coordinates is a metric function, strings can be compared by edit distance, the document full-text search uses cosine distance, sets can be compared by Jaccard coefficient or Hausdorff distance, vectors can use L_p metrics or Earth mover's distance, and so on. From that point of view, the MUFIN system is highly a *extensible* multi-purpose engine that can supply efficient similarity search in various areas.

The MUFIN system achieves the search *scalability* by adopting paradigms of structured peer-to-peer networks: (1) dynamically distributing workload on independent peers for potential parallel query execution, (2) avoiding bottlenecks formed by a single entry point or centralized directories - typical for traditional client-server architectures, and (3) improving fault tolerance by replication of data and adopting multiple search strategies. It is possible by using the concept of virtual processing units (peers) to which the MUFIN system maps the respective parts of the search engine. Due to this open architecture, the system can adapt to virtually any amount of data and also the throughput of the system – the number of simultaneous queries that can be solved without slowing down – is increasing, since queries can be posed from any peer. Additionally, peers can be replicated transparently on the system level not only improving the system robustness but also its throughput. The concept also allows running several indices for different aspects of the data, e.g. facial features, fingerprints and palmprints of a person, together. The system then makes a combined search possible, e.g. it can identify people according to several methods while taking the respective methods' accuracies into account.

Even though the ideas of the system design come from structured peer-to-peer networks, MUFIN is perfectly suitable to run in more controlled environments like computer clusters, GRIDs or even multi-CPU servers. Actually, implementing a MUFIN searching service on computing clouds (e.g. Amazon EC2) yields a cost-effective fully scalable solution with guaranteed availability. The system hardware abstraction layer allows MUFIN to *map* its peers onto different hardware architectures. This allows tuning the performance of the system dynamically – if a higher throughput or faster query response is required the system can utilize additional hardware without changes in the underlying indices. On the other hand, when the system is under-loaded (e.g. at night or during weekends) it can shrink freeing the computational resources that can be used by other tasks. The MUFIN system also maintains the actual object data independently from the index structures thus allowing effective utilization of various storage architectures as well as caching the data in main memory.

3.0 MUFIN IMPLEMENTATION

The implementation of the MUFIN system builds on a framework called Metric Similarity Search Implementation Framework (MESSIF) [1] which is designed to ease the task of building metric-based indexing techniques. Basically, the framework is a collection of modules schematically depicted in Figure 2. The metric space module encapsulates support for metric data objects, the operations module offers a framework for data querying and manipulation methods, and the storage module provides interfaces for creating and maintaining various data storages. The communication module allows exchanging information via the computer network by means of sending and receiving messages. Modules for centralized and distributed index structures encapsulate implemented indexing techniques providing a unified access to searching and data manipulation regardless of the particular implementation details. Finally, the statistical module allows gathering performance characteristics of all the other modules, e.g. it can be used to monitor the number of I/O operations during the search or to measure the amount of data sent over the network.

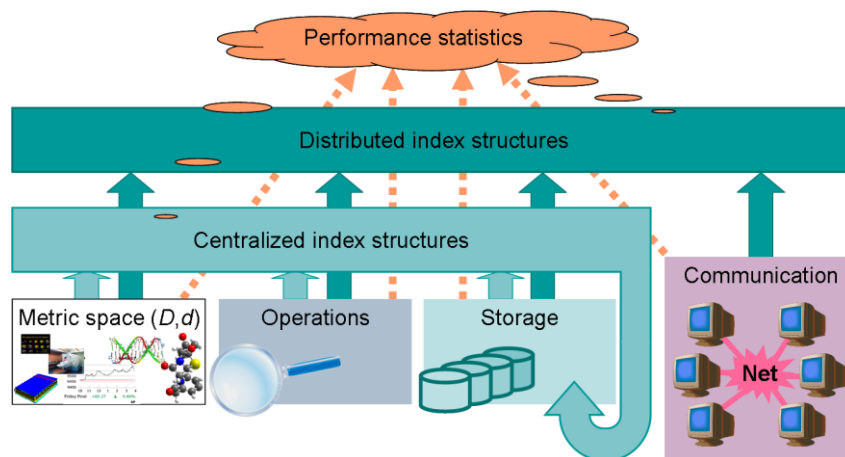


Figure 2: Overview of Metric Similarity Search Implementation Framework (MESSIF).

A wide variety of metric indexing techniques is integrated into the MUFIN system – from traditional tree-like structures (e.g. M-Tree [2], PM-Tree) through governor-based parallel indexes (e.g. M-Grid [3]) to fully distributed indices (e.g. GHT* [4], M-Chord [5]). Since the system has a unified public interface defined for the indices, integrating additional techniques or linking with external services is easy. Data within the system can be manipulated and queried via several user interfaces. The batch interface is suitable for building huge indices or insertions of data bulks as well as running automatic experiments. Developer APIs are available for linking with other systems using standardized technologies like WebServices or Java RMI. And finally, specialized web application interfaces allow users to interact with the system web browser.

4.0 DEMONSTRATION APPLICATIONS

In this section, we illustrate capabilities of MUFIN by several applications which use various data types with different ways of assessing similarity. Presented systems use datasets of diverse volumes and run on various hardware infrastructures.

4.1 Large-scale General Image Search

The first application is a large-scale searching system for general images from the World Wide Web. The images are searched according to their content. It demonstrates both the applicability of MUFIN to this issue and its ability to efficiently manage Web-scale datasets [6].

Data

The dataset consists of 100 million digital images from the CoPhIR dataset¹. The images were crawled from a photo-sharing system Flickr² and preprocessed – five different visual descriptors were extracted from each of the image. These descriptors are specified in MPEG-7 standard [8] and they capture various color, shape and texture characteristics of the image. For each of the descriptors, there is a metric function to measure similarity of two images with respect to the descriptor [9]. We combine the five descriptors into a single metric space by means of a weighted sum of individual descriptors' distances. See Table 1 for details about the used descriptors and their weights.

¹ CoPhIR: Content-based Photo Image Retrieval Collection [7]: <http://cophir.isti.cnr.it>

² <http://www.flickr.com>

MPEG-7 Descriptor	Metric	Weight
Scalable Color	L1 metric	2
Color Structure	L1 metric	3
Color Layout	sum of L2 metrics	2
Edge Histogram	weighted sum of L1 metrics	4
Homogeneous Texture	weighted sum of L1 metrics	0.5

Table 1: CoPhIR visual descriptors and their combination.

Index and Infrastructure

This specific instance of the MUFIN searching system is formed by a distributed index called M-Index [10] that is based on peer-to-peer principles. This architecture is highly scalable and also flexible in terms of the hardware infrastructure it runs on – in this case, the index is formed by 500 logical peers hosted by six IBM servers (8 CPU cores, 16GB RAM, and 6 disks). The performance of the system (in terms of response time and query throughput) can be directly tuned by altering the HW infrastructure.

Example

In general, large volumes of data have positive influence on quality of similarity search results. Figure 3 demonstrates this “power of volume” – the upper and lower rows show results of MUFIN similarity search in 10 and 100 million Flickr images, respectively (the first image is the “query” and the rest are the most similar images). We can see a noticeable improvement in the search effectiveness caused by growing number of images that are more similar to the query image. This application is publically available at <http://mufin.fi.muni.cz/imgsearch/>.



Figure 3: Results of similarity search in 10 million and 100 million Flickr images.

4.2 Multiple Visual Aspects

In the system described in the previous section, the images are searched according to a single combination of various visual aspects – colors, shapes, texture. In the following application, we have created multiple indexes to let user select which visual aspect to use for searching.

Data

In this application, we use a dataset of 10,000 images of e-shop commodities. The following MPEG-7 descriptors were extracted from every image: Region Shape, Edge Histogram, Scalable Color, Color Structure, and Color Layout. These descriptors can be now combined in a specific way to emphasize specific visual aspects.

Index and Infrastructure

We have built separate indexes for the following three combinations of visual descriptors:

- shapes: Edge Histogram and Region Shape;
- colors: Scalable Color, Color Structure, and Color Layout;
- shapes & colors: combination of all five descriptors.

Because each of these combinations forms a metric space, they can be straightforwardly managed by MUFIN. Namely, all the indexes use the M-Index structure and can be stored in main memory because the dataset is not large.

Example

Figure 4 shows an example of similarity search for one query image according to the three combinations. The first row shows result images searched with a focus on shapes, the second on colors, and the third with contribution of both aspects.



Figure 4: Example of similarity search by multiple visual aspects.

4.3 Local Descriptors

Both the applications introduced above use *global visual descriptors*, which capture visual characteristics from the image as from a whole. The *local descriptors* [11] such as SIFT [12] work in a slightly different way – first an algorithm identifies so-called *keypoints* in the images and then special descriptors are extracted from each of the keypoint.

Data and Index

The dataset is formed by approximately 15,000 logo images downloaded from the Web. From each image, we extract both the five MPEG-7 global descriptors described in the previous section and SURF local descriptors [13]. The global descriptors are combined and managed in the same way as in the applications above and the sets of local descriptors extracted from each image are compared by a special metric function and managed by a disk-oriented version of the M-Index structure.

Example

Figure 5 shows an example of search results using the global and local visual descriptors. We can see that while the global descriptors capture the overall composition of the image, local descriptors can identify relations between fragments of the images, for instance a similar font used in various parts of the images.



Figure 5: Results of similarity search using global and local visual descriptors on a logo dataset.

4.4 Searching by GPS Location

It becomes more and more common that digital photographs contain information about the GPS location where the image was taken. This piece of information is either created automatically by the digital camera itself or added later by user, and it is typically stored within *Exif* metadata (Exchangeable image file format). The GPS locations can be used for searching.

Data and Index

The 100 million CoPhIR dataset, described in the first application section, is composed mainly by digital photos and about 9% of them contain GPS location information. We use this 9 million dataset and search it according to the space distance between individual locations. As this distance is naturally a metric we can use standard MUFIN technologies to implement this index and search.

Example

Figure 6 shows an example of the MUFIN GPS search. User specifies GPS coordinates by selecting a location at the map and the results are photos taken near the specified location. This search system was created within European project SAPIR, and MUFIN is used as a service called from the user frontend.



Figure 6: MUFIN image search using the GPS locations stored in digital photographs.

4.5 Video Similarity Search

Recently, we can observe a mass production of audiovisual data. As well as other complex data types, video can be either searched by annotations or by content of the data itself. As MUFIN is based on general principles, it provides a number of ways to index and search such data by their content.

Data and Index

In the video application, we use a sample of 40GB of news video (15 hours) provided by BBC for European project SAPIR. The data are pre-processed in the following way: standard software is used to identify *keyframes* in the video sequences and then five MPEG-7 descriptors are extracted from each keyframe in the same way as in the application introduced in Section 4.1. These descriptors are indexed and searched as digital images.

Example

Figure 7 shows an example of the result from the BBC news video search. Please, note that global visual descriptors were used in this application to capture the overall composition of the scene – in this example it is “a newscaster in a studio with the same background image”.



Figure 7: Example of audiovisual data search according to keyframe similarity.

5.0 THE CHALLENGE OF BIOMETRIC DATA

In general, biometrics is automated methods of recognizing a person based on the person’s physiological and behavioral characteristics. Biometrics include a wide variety of technologies ranging from traditional fingerprints over facial or iris recognition and retinal scanning to DNA testing, speech verification and gait recognition. There are two types of biometric recognition problems:

- **Verification** – the aim is to verify whether the person is who he/she claims to be, i.e. verifying authenticity of the person;
- **Identification** – a person exposes biometric characteristic and the aim is to tell who the person is.

The main advantage of most biometric characteristics is their uniqueness and stability over time, i.e. they usually do not change over a short period of time. Biometric characteristics cannot be lost, forgotten or stolen, so it prevents impersonation and repudiation. It is also very hard, or at least unusual, to hide the characteristics, e.g. gait (the way a person walks) is very difficult to change.

Accuracy of biometric devices depends on the quality of the scan of a biometric characteristic. This leads to *false rejects* or *false accepts*. These errors can be handled by improving capturing capabilities of individual biometric devices or by using a more different characteristic and proceeding *biometric data fusion*. There are several levels of biometric data fusion:

- **Sensor level** – a person’s characteristic is captured multiple times and the sensor/device combines the scans to improve the resulting quality or multiple sensor are used to capture data;
- **Feature level** – several algorithms can be used to extract different features from the same biometric scan, i.e. each of the algorithms can focus on a distinct property of the scan;
- **Instance level** – more instances of the same kind are used to capture biometric data, e.g. fingerprints of more fingers or iris scans of both eyes;
- **Characteristic level** – different biometric characteristics are used to verify/identify a person, e.g. face and fingerprint.

A solution to the problem of identification clearly requires a similarity search engine capable of online processing large databases. In verification, such an engine does not necessarily be applied since a person states its identity and the person’s authenticity is verified. In particular, biometric scans can be retrieved from a database by the person’s ID, so a conventional primary-key search technique can be used. However, the application may require testing biometric scans for uniqueness during a new person’s enrolment to the system, i.e. the new person is added to the database. In this case, a similarity search engine is a must. On the other hand, applying a biometric technology usually requires significant computational resources. The MUFIN technology is directly applicable to these aforementioned issues since it offers similarity search and scalability in terms of computational and storage resources while retaining an acceptable query response time. In the following, we summarize details about selected biometric modalities.

5.1 Fingerprints

Famous applications of fingerprints are in criminalistics and in border and immigration control. Minutiae is one successfully applied method of comparing ridges in fingerprints [14]. It identifies places where ridges start, stop or bifurcate (branch), refer to Figure 8. These places are then observed as points with a direction and are converted to polar coordinates. As a result, a fingerprint is described as a sequence of points in polar coordinates. Two sequences are then matched using a weighted edit distance function. The weights used do not break metric postulates, so this distance function is directly applicable to MUFIN.



Figure 8: Example of minutiae extracted from a fingerprint image (taken from [14]).

Another approach to match fingerprints uses both local and global characteristics of each fingerprint, which is the advantage over the previous technique that focuses on local ones only. This approach is called Filterbank-based matching [15]. Each fingerprint is divided to regions and Gabor filters detecting a certain direction of ridges are applied. To catch global and local characteristics, the filters are tuned to eight

directions (by 22.5 deg.), see Figure 9. The individual regions are then expressed as a single number (average absolute deviation from the mean of gray pixel values). These values form a high-dimensional vector. The Euclidean distance is applied to measure the similarity.

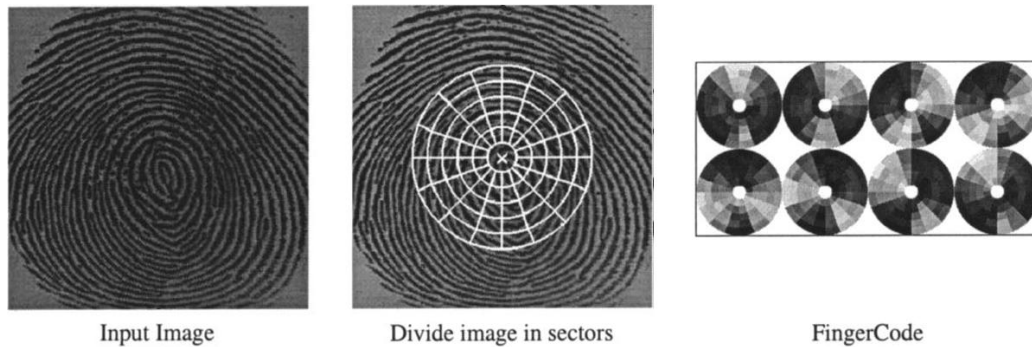


Figure 9: Extraction of Filterbank-based features (taken from [15]).

5.2 Hand recognition

Hand biometrics are successfully commercially used, e.g. to verify identity of employees in nuclear power plants. Basically, there are three types of biometric features extracted from a hand image: hand geometry, palmprint and finger surface. In [17], the authors propose an algorithm that extracts hand silhouettes from regular scans. After removing the background, a hand contour is obtained. It is further processed by global registration that rotates and translates the contour to a standard position. Next, ring artefacts are removed and individual fingers are rotated, please refer to Figure 10. Contours of two hands are compared by a modified Hausdorff distance. The modifications made by the authors do not break the metric function properties, thus it is directly applicable to MUFIN.

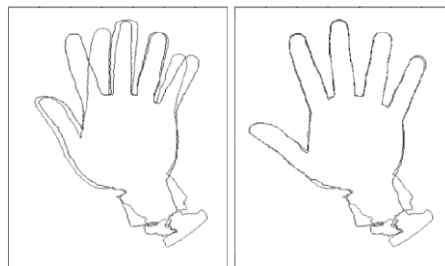


Figure 10: Example of registering two hand contours of the same person: (left) fingers are not registered, (right) fingers are registered (taken from [17]).

5.3 Gait Recognition

Gait, or the way a person walks, is a unique and idiosyncratic characteristic of the person. Its advantage for biometrics is that it is difficult to conceal and it can be easily captured even at long distances. In paper [16], the gait information is extracted from a video sequence. In particular, a silhouette of the walking person is determined for each video frame by subtracting the background of the image. The sequence of silhouettes is divided in subsequence, each of them representing one gait cycle (two steps). Then, an average silhouette is computed for each subsequence, please refer to Figure 11. The binary silhouettes are then compared using the Euclidean distance.

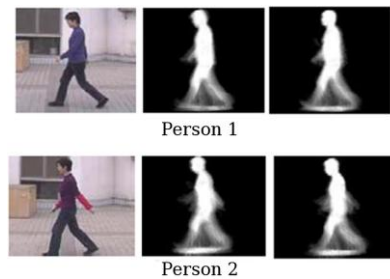


Figure 11: Example of average silhouette extraction (taken from [16]).

5.4 Face Recognition

Photographs of faces are widely used, for example, in ID documents. An advantage is that few people object to having their photo taken, so face recognition as a biometrics discipline can be used. Another advantage, shared with the gait recognition, is that photographs can be taken covertly and at long distances. On the other hand, the precision of face recognition can be influenced by physiologic changes as such growing facial hair. Basically, there are two approaches to process faces: image-based and feature-based. The former analyses the raster image of the face while the latter captures geometric characteristics or other metrics such as spatial relationships. In the following, we describe an application of MUFIN that uses the image-based face recognition.

Data and Index

In this system, we use a collection of 10,000 photos from a photo database of Masaryk University. We have used publicly available software to detect 16,000 faces in these photos and then extracted MPEG-7 *Advanced Face Descriptor* [9] for these faces. This descriptor operates on the normalized image. By applying principle component analysis (PCA), eigenvectors are extracted. The eigenvectors with large eigenvalues capture information that is common for a group of faces while the vectors with small values represent information specific to a particular face. The similarity of descriptor vectors is evaluated by the Euclidean distance, which again allows the straightforward application to MUFIN.

Example

Figure 12 shows an example of result from the MUFIN face retrieval application. We can see a query face (in the red rectangle) and the search results. The seven found faces demonstrate that the query face has been found in the Masaryk University database and no false faces have been returned. On the right, the original photo with detected faces for one of the retrieved faces is shown.



Figure 12: Face recognition application of MUFIN.

6.0 CONCLUSIONS

There are no doubts that the modern similarity search paradigm finds also a lot of applications in terrorist informatics. In this paper, we have briefly introduced the MUFIN system which serves as a universal engine for similarity searching. The MUFIN architecture is designed on the concepts of: (1) extensibility – to achieve applicability for diverse data collections comparing objects by varied measures of similarity, (2) scalability – to also process extremely large collections of data queried by many concurrent requests, and (3) infrastructure independence – to tune performance according to requirements of specific applications. We have also exemplified the capabilities of MUFIN by several demonstration applications. Then we have analyzed some biometric descriptors and shown that they can be easily processed by MUFIN, because most of them satisfy the metric properties. Finally, we have demonstrated the universality of MUFIN by the face retrieval application which looks for the most similar faces stored within a database to a given one. Our future research directions focus on large-scale experiments on varied biometric characteristics.

ACKNOWLEDGEMENTS

This work has been supported by national projects GA201/09/0683, GP201/08/P507, GP201/07/P240, GD102/09/H042, and MUNI/E/0066/2009. Hardware infrastructure has been provided by MetaCenter (<http://meta.cesnet.cz>).

REFERENCES

- [1] Michal Batko, David Novák and Pavel Zezula. MESSIF: Metric Similarity Search Implementation Framework. In *Digital Libraries: Research and Development, Lecture Notes in Computer Science*, Berlin, Heidelberg : Springer-Verlag, vol. 4877, 10 pages, 2007. ISBN 978-3-540-77087-9.
- [2] Paolo Ciaccia, Marco Patella and Pavel Zezula. M-tree: An Efficient Access Method for Similarity Search in Metric Spaces. In *Proceedings of the 23rd International Conference on Very Large Data Base*, San Francisco, CA, USA, Morgan Kaufmann Publishers Inc. 10 pages, 1997. ISBN 1-55860-470-7.
- [3] Michal Batko, Vlastislav Dohnal and Pavel Zezula. M-Grid: Similarity Searching in Grids. In *Proceedings of International Workshop on Information Retrieval in Peer-to-Peer Networks, ACM CIKM 2006*. Arlington, ACM Press. 8 pages, 2006. ISBN 1-59593-531-2.
- [4] Michal Batko, Claudio Gennaro and Pavel Zezula. Scalable and Distributed Similarity Search in Metric Spaces. In *5th Workshop on Distributed Data & Structures. Revised Selected Papers*. Canada, Carleton Scientific. 10 pages, 2004. ISBN 1-894145-18-6.
- [5] David Novák and Pavel Zezula. M-Chord: A Scalable Distributed Similarity Search Structure. In *InfoScale '06: Proceedings of the 1st international conference on Scalable information systems*. New York, NY, USA, ACM Press. 10 pages, 2006. ISBN 1-59593-428-6.
- [6] David Novak, Michal Batko and Pavel Zezula. Web-scale System for Image Similarity Search: When the Dreams Are Coming True. In *Proceedings of the Sixth International Workshop on Content-Based Multimedia Indexing (CBMI 2008)*, London, UK, IEEE Computer Society, Los Alamitos, CA 90720-1314. 8 pages, 2008.

- [7] Paolo Bolettieri, Andrea Esuli, Fabrizio Falchi, Claudio Lucchese, Raffaele Perego, Tommaso Piccioli and Fausto Rabitti. CoPhIR: a Test Collection for Content-Based Image Retrieval. *CoRR Journal*, vol. abs/0905.4627v2. 15 pages, 2009.
- [8] MPEG-7. Multimedia content description interfaces. Part 3: Visual. ISO/IEC 15938-3:2002. 2002.
- [9] B.S. Manjunath, Phillipe Salembier and Thomas Sikora. Introduction to MPEG-7: Multimedia Content Description Interface. John Wiley & Sons, Inc., New York, NY, USA. 396 pages, hardcover, 2002. ISBN 978-0-471-48678-7.
- [10] David Novak and Michal Batko. Metric Index: An Efficient and Scalable Solution for Similarity Search. In *2nd International Workshop on Similarity Search and Applications*, Prague, Czech Republic, IEEE Computer Society, Los Alamitos, CA 90720-1314. 9 pages, 2009. ISBN 978-0-7695-3765-8.
- [11] K. Mikolajczyk, T. Tuytelaars, C. Schmid, A. Zisserman, J. Matas, F. Schaffalitzky, T. Kadir and L. Van Gool. A Comparison of Affine Region Detectors. *International Journal of Computer Vision*, Kluwer Academic Publishers, Hingham, MA, USA. vol. 65, pages 30, 2005.
- [12] David G. Lowe. Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision*. vol. 60, pages 20, 2004.
- [13] Herbert Bay, Tinne Tuytelaars and Luc Van Gool. SURF: Speeded Up Robust Features. *Lecture Notes in Computer Science*. vol. 3951, 14 pages, 2006.
- [14] Davide Maltoni, Dario Maio, Anil K. Jain and Salil Prabhakar. Handbook of Fingerprint Recognition. Springer. 496 pages, hardcover, 2nd edition, 2009. ISBN 978-1-84882-253-5.
- [15] Anil K. Jain, Salil Prabhakar, Lin Hong and Sharath Pankanti. Filterbank-Based Fingerprint Matching. *IEEE Transactions on Image Processing*. vol. 9, pages 14, 2000.
- [16] J. Fazenda, D. Santos and P. Correia. Using Gait to Recognize People. In *Proceedings of the International Conference on Computer as a Tool (EUROCON 2005)*, pages 4, 2005. ISBN 1-4244-0049-X.
- [17] E. Yoruk, E. Konukoglu, B. Sankur and J. Darbon. Shape-based hand recognition. *IEEE Transactions on Image Processing*. vol. 15(7), pages 13, 2006.

